

§ Bandits with general function approximation

Setup: For $t=1, 2, \dots, T$:

observe context $x_t \in \mathcal{X}$

take action $a_t \in A$

receive reward $r_t = f^*(x_t, a_t) + \underline{\epsilon}_t$ 1-SG zero mean

Goal: minimize Pseudo regret

$$P\text{Reg}(T) = \sum_{t=1}^T \max_a f^*(x_t, a) - \sum_{t=1}^T f^*(x_t, a_t) \quad \text{① } F \subseteq (\mathcal{X} \times A \rightarrow [0])$$

Assumption: $f^* \in \bar{F}$ (realizability) ② $|F| < \infty$.

But functions in \bar{F} are not nec. $f(x, a) \neq \langle \theta, \phi(x, a) \rangle$
in case

Example: generalized linear bandits: linear bandits w/ an "activation fn".

$$f^* \in \bar{F} = \left\{ f_\theta(x, a) = \sigma(\langle \theta, \phi(x, a) \rangle) : \|\theta\|_2 \leq 1 \right\}$$

$$\textcircled{1} \quad r_t \mid x_t = x, a_t = a \sim \text{Bernoulli}\left(\frac{1}{1 + e^{-\langle \theta^*, \phi(x, a) \rangle}}\right)$$

logistic bandits (suitable for modeling binary responses - e.g. clicks)

$$\Rightarrow f^*(x, a) = \frac{1}{1 + e^{-\langle \theta^*, \phi(x, a) \rangle}}$$

$$\textcircled{2} \quad r_t \mid x_t = x, a_t = a \sim \text{Poisson}\left(e^{\langle \theta^*, \phi(x, a) \rangle}\right)$$

Poisson rewards (modeling counts, e.g. number of visits)

$$\Rightarrow f^*(x, a) = e^{\langle \theta^*, \phi(x, a) \rangle}$$

$\approx \text{Bin}(n, \frac{\lambda}{n})$ for large n .

- $\text{Poisson}(X; \lambda) = \frac{\lambda^x}{x!} e^{-\lambda} \quad x \in \mathbb{N}$
 $\mathbb{E}[X] = \lambda$
- ③ $r_t | x_t = x, a_t = a \sim N(\langle \theta^*, \phi(x, a) \rangle, 1) \Rightarrow f^*(x, a) = \langle \theta^*, \phi(x, a) \rangle$
- ④ other examples: e.g. Exponential distn,

Gamma. Pareto. Weibull (key word: exponential family)
 (Filippi, Cappé, Garivier, Szepesvári, 2010).

Given a nonlinear F . how to design low-regret algorithm?

Again: optimism principle.

Algorithm (OFU : optimism in the face of uncertainty)

World model: reward model f

For $t=1, 2, \dots T$:

- construct confidence set F_t for f^* .
 - Observe x_t .
 - Take action $a_t = \underset{a \in A}{\operatorname{argmax}} \underset{f \in F_t}{\max} f(x_t, a)$
- The largest plausible reward action a can get

Again. Q1 How to construct tight F_t . $\forall t$?

Q2 How to relate the tightness of F_t to alg's regret?

for Q1: can mimic linear bandits.

- ① define a center \hat{f}_t ② let F_t be a ball around

at \hat{f}_t with certain distance metric.

① Natural idea: "best fit".

$$\hat{f}_t = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{s=1}^{t-1} (f(x_s, a_s) - r_s)^2 =: L_{t-1}(f)$$

② How to define distance metric?

Recall: in linear bandits, cannot guarantee \hat{f}_t and f^* close everywhere.

Measure closeness using "data norm".

$$S_{t-1} = \{(x_s, a_s)\}_{s=1}^{t-1}$$

$$\mathcal{F}_t = \{f \in \mathcal{F}: \|f - \hat{f}_t\|_{S_{t-1}}^2 \leq \beta_t\}$$

Here $\|g\|_S = \sqrt{\sum_{(x,a) \in S} g(x,a)^2}$ is the "data norm" of g wrt dataset S .

$$\text{so } \mathcal{F}_t = \{f \in \mathcal{F}: \sum_{s=1}^{t-1} (f(x_s, a_s) - \hat{f}_t(x_s, a_s))^2 \leq \beta_t\}$$

shorthand: $\mathbb{Z}_t := (x_t, a_t) \cdot \forall t$.

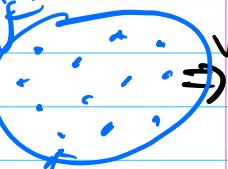
Sanity check: for linear bandits. if $f \leftrightarrow \theta$.
 $\hat{f}_t \leftrightarrow \hat{\theta}_t$

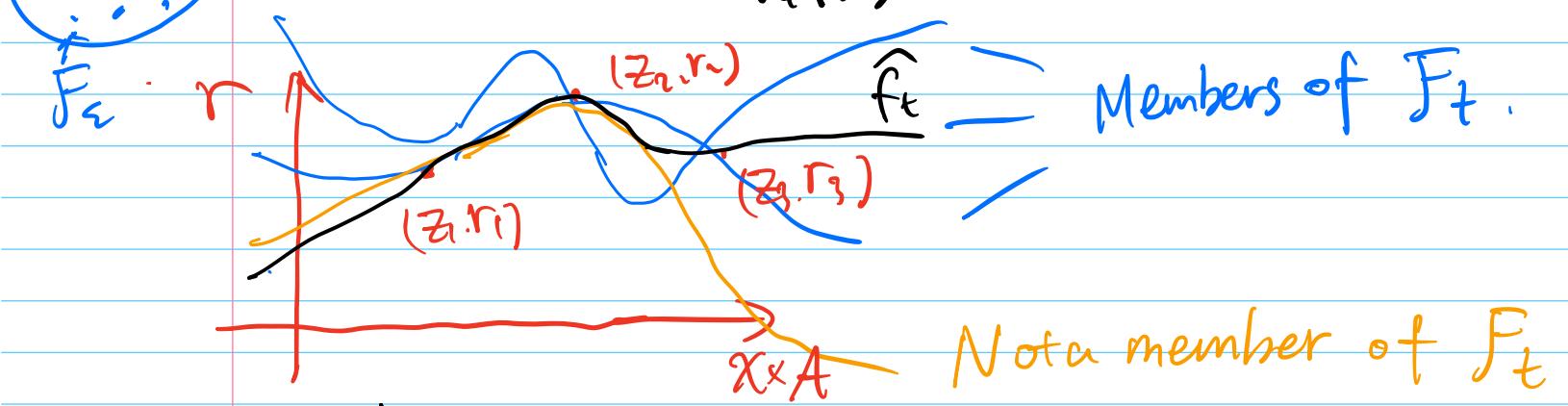
$$\text{The above } \Rightarrow \sum_{s=1}^{t-1} (\langle \theta - \hat{\theta}_t, \phi(x_s, a_s) \rangle)^2 \leq \beta_t$$

$$\Leftrightarrow \sum_{s=1}^{t-1} (\theta - \hat{\theta}_t)^\top \phi_s \phi_s^\top (\theta - \hat{\theta}_t) \leq \beta_t$$

$$\Leftrightarrow \|\theta - \hat{\theta}_t\|_{V_{t+1}(0)}^2 \leq \beta_t$$

Lemma (Validity of conf. set) Let $\beta_t = 8 \ln \frac{4|F|}{\delta}$ $\forall t$, Then, w.p. $1-\delta$:
 $\forall t \quad f^* \in F_t$ (Validity of confidence set)

(Recall: for linear bandits, $\ln |F| \geq \# \text{parameters} = d$.

For infinite F , can consider its discretization replaced by its covering number.
 $\Rightarrow \forall t \quad \|\theta - \hat{\theta}_t\|_{V_{t+1}(0)}^2 \leq \tilde{\Omega}(d)$ same as before)



Pf of validity lemma:

Claim: there exists event E , $P(E) \geq 1-\delta$.

on E : $\forall f \in F$, $\forall t$,

$$\frac{1}{2} \|f - f^*\|_{S_t}^2 - 4 \ln \frac{2|F|T}{\delta} \stackrel{\textcircled{1}}{\leq} L_t(f) - L_t(f^*) \stackrel{\textcircled{2}}{\leq} \frac{3}{2} \|f - f^*\|_{S_t}^2 + 4 \ln \frac{4|F|T}{\delta}$$

Interpretation: ① If $L_t(f)$ small $\Rightarrow f$ close to f^* .

② If f is close to f^* $\Rightarrow L_t(f)$ small.

Claim \Rightarrow lemma: on E , want to show $\|\hat{f}_t - f^*\|_{S_t}^2$

small. ($\leq \beta$)

use ① or ②? ① ✓

$$\Rightarrow \frac{1}{2} \| \hat{f}_t - f^* \|_{S_t}^2 - 4 \ln \frac{2\text{FI}T}{\delta} \leq L_t(\hat{f}_t) - L_t(f^*)$$

optimality of \hat{f}_t
 ≤ 0

$$\Rightarrow \| \hat{f}_t - f^* \|_{S_t}^2 \leq 8 \ln \frac{2\text{FI}T}{\delta} = \beta. \quad \square.$$

proof of the claim :

It suffices to show $\forall f, \forall t$, w.p. $1-\delta$:

$$\frac{1}{2} \| f - f^* \|_{S_t}^2 \stackrel{\textcircled{a}}{\leq} L_t(f) - L_t(f^*) \stackrel{\textcircled{b}}{\leq} \frac{3}{2} \| f - f^* \|_{S_t}^2 + 4 \ln \frac{2}{\delta}$$

~~$- 4 \ln \frac{2}{\delta}$~~

(Why? Define $E = \bigcap_{f \in F} \bigcap_{t=1}^T E_{f,t}$ and apply union bounded (exercise))

we will prove ② (① is similar)

Step 1: understand what $L_t(f) - L_t(f^*)$ concentrates to; (a r.v. should concentrate to its expectation)

$$L_t(f) - L_t(f^*) = \sum_{s=1}^t (f(z_s) - r_s)^2 - (f^*(z_s) - r_s)^2$$

If $z_1 \dots z_t$ are all chosen ahead of time & deterministically, what's its expectation?

$$= \sum_{s=1}^t (f(z_s) - f^*(z_s) - \varepsilon_s)^2 - \varepsilon_s^2$$

The expectation here is w.r.t to the randomness of noise $\varepsilon_1 \dots \varepsilon_t$.

$$= \sum_{s=1}^t (f(z_s) - f^*(z_s))^2 - \sum_{s=1}^t 2(f(z_s) - f^*(z_s)) \varepsilon_s$$

$\underbrace{\|f-f^*\|_{S_t}^2}$, fixed. $\underbrace{2(f(z_s) - f^*(z_s)) \varepsilon_s}$, expectation zero.

Therefore, we focus on understanding the concentration of

$$L_t(f) - L_t(f^*) = \|f - f^*\|_{S_t}^2$$

$$= - \sum_{s=1}^t 2(f(z_s) - f^*(z_s)) \varepsilon_s =: M_t$$

$\{M_t\}$ is a martingale & $E[M_t] = 0 \cdot \forall t$.

At this point, can use e.g. Azuma's inequality

or self-normalized tail ineq. for martingales

(linear bandit lecture) to bound the RHS.

However: ① this may give suboptimal bounds,

② there is a more direct way:

Step 2: utilize "subgaussian" property of M_t

Intuition: If all $z_1 \dots z_t$ are chosen ahead,

M_t is SG with variance proxy $4\|f - f^*\|_{S_t}^2$

i.e. $\mathbb{E} e^{\lambda M_t - 2\lambda^2 \|f - f^*\|_{S_t}^2} \leq 1 \quad \forall \lambda \in \mathbb{R}$

Fact: (*) continues to be true when $z_1 \dots z_t$

are chosen adaptively online.

(proof: use law of iterated expectations)

Fix $\lambda > 0, \gamma > 0$.

Markov's ineq $\Rightarrow \lambda M_t - 2\lambda^2 \|f - f^*\|_{S_t}^2 \leq \ln \frac{1}{\gamma}$

$$\Rightarrow L_t(f) - L_t(f^*) \leq ((+2\lambda)) \|f - f^*\|_{S_t}^2 + \frac{1}{\lambda} \ln \frac{1}{\gamma}$$

take $\lambda = \frac{1}{4}$: $\cancel{\lambda}$

Q2: Starting point:

Note: under E. For OFU

$$PReg(T) = \sum_{t=1}^T \max_{a \in A} f^*(x_t, a) - f^*(x_t, a_t)$$

same analysis as in OFUL

$$\leq \sum_{t=1}^T \max_{f \in F_t} f(x_t, a_t) - f^*(x_t, a_t)$$

$$f^+ \in F_t \leq \sum_{t=1}^T \max_{f \in F_t} f(x_t, a_t) - \min_{f \in F_t} f(x_t, a_t)$$

!! "width/uncertainty" of z_t wrt historical observations.

key insight: in MAB and linear bandits. $w_{F_t}(z_t)$ do not stay high for a long time. due to interdependences of the z_t 's seen.

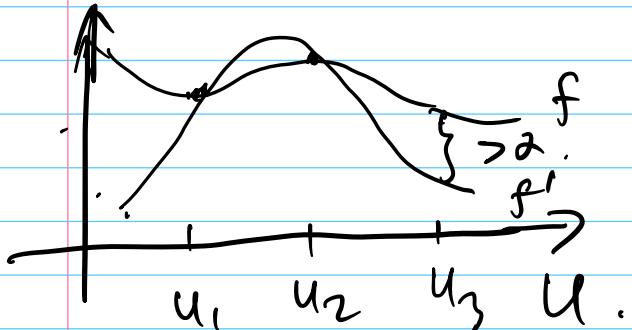
Can we characterize the interdependence of $\{z_t\}$'s for general F ?

Def (Eluder dimension, Russo & Van Roy, 2013).

- u_i is α -independent of $\underbrace{u_1 \dots u_{i-1}}_{U_{i-1}}$. with respect

to F . if

$$\exists f, f' \in F. \|f - f'\|_{U_{i-1}}^2 \leq \alpha^2 \text{ but } |f(u_i) - f'(u_i)| > \alpha.$$



(otherwise. u_i is said to be α -dependent of U_{i-1})

i.e. knowing f^* in U_{i+1} will imply knowledge of approximate $f^*(u_i)$.

The Eluder dimension of F at scale ε , $\text{Edim}(F, \varepsilon)$

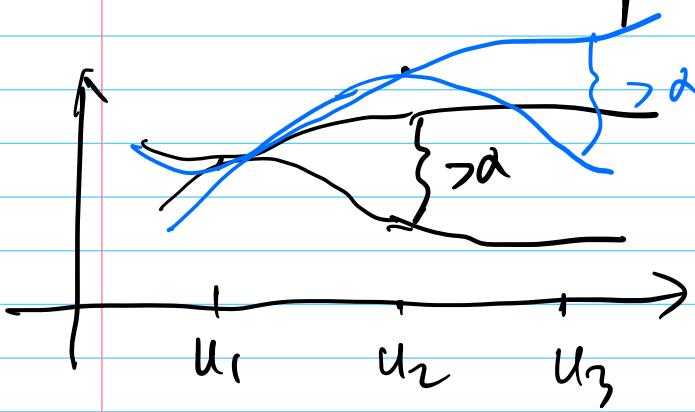
is the length of the longest sequence $u_1 \dots u_N$.

so that $\exists \alpha > \varepsilon$, and:

- u_2 is α -indep of u_1
- u_3 is α -indep of u_1, u_2

...

- u_N is α -indep of $u_1 \dots u_{N-1}$



... (analogy: $u_1 \dots u_N$ is some "independent basis", similar to linear independence in linear algebra.)

Politician tries to evade reporters:

(f^*)

- Hope to keep her true position hidden
- Each piece of info provides needs to be new
(cannot be inferred from previous ones) $(f^*(u_i))$
- How long can she continue before her

$\leftarrow X^* A$
|
| $u_1 \dots u_N$
| on α -
| indept
| sequence

position is determined?

Examples of Edim:

- $\text{Edim}(F, \epsilon) \leq |\mathcal{X} \times A|$

(an α -indep sequence cannot have repetitions)

- $F = \{ \sigma(\langle \theta, \phi(x, a) \rangle) : \| \theta \|_2 \leq 1, \underbrace{0 < L_- \leq \sigma'(z) \leq L_+ \wedge z}$

$$\text{Edim}(F, \epsilon) \leq \tilde{\mathcal{O}}\left(\left(\frac{L_+}{L_-}\right)^2 \cdot d \ln \frac{1}{\epsilon}\right)$$

"well-behaved activation fn".

- (Li, Kamath, Foster, Srebro: "understanding the Eudler Dimension", 2022)

Regret analysis of OFU

Thm: w.p. $1-\delta$: $\text{PReg}(T) \leq \tilde{O}(\sqrt{\text{Edim}(F, \frac{1}{T})} \cdot \ln |F| T)$

For generalized linear bandits:

$$\text{Edim}(F, \frac{1}{T}) = \tilde{\mathcal{O}}\left(\left(\frac{L_+}{L_-}\right)^2 d \ln T\right) \quad \ln |F| = \tilde{\mathcal{O}}(d)$$

$$\Rightarrow \text{PReg}(T) \leq \tilde{\mathcal{O}}\left(\left(\frac{L_+}{L_-}\right)^2 d \sqrt{T}\right).$$

linear bandits: $L_+ = L_- = 1$.

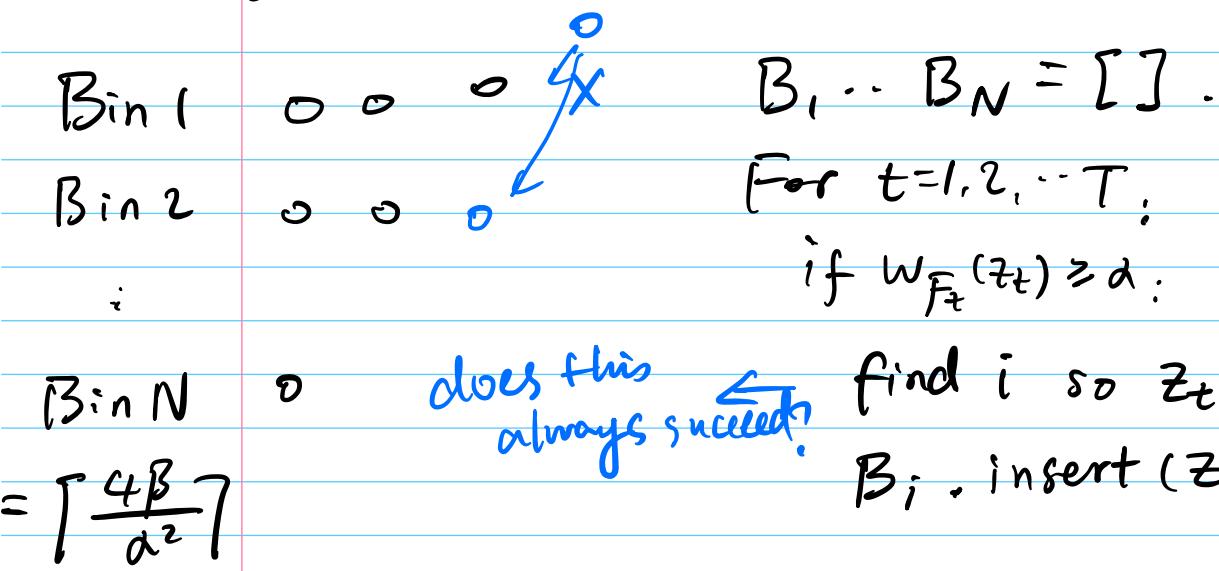
Proof: Claim: $\forall \alpha > 2$,

$$\sum_{t=1}^T I(W_{F_t}(z_t) > \alpha) \leq \left(\frac{4\beta}{\alpha^2} + 1\right) \text{Edim}(F, \Sigma).$$

Intuition: #rounds of "surprises" should be small.
 \downarrow
(large $w_{F_t}(z_t)$)

so the total surprises should also be small.

Pf of Claim: we will arrange $\{z_t : w_{F_t}(z_t) \geq \alpha\}_{t=1}^T$ in a way so that we can see its size is controlled.



Observation 1: the "Find" step always succeeds.

Reason: let \bar{f} and f be in F_t s.t.

$$\bar{f}(z_t) - f(z_t) = w_{F_t}(z_t) > \alpha$$

$$\text{Note: } \| \bar{f} - \hat{f}_t \|_{S_{++}}^2 \leq \beta$$

$$\|\underline{f} - \hat{f}_t\|_{S_{t-1}}^2 \leq \beta$$

$$\Rightarrow \|\bar{f} - \underline{f}\|_{S_{t-1}}^2 \leq 4\beta \quad (\text{exercise}) \quad (*)$$

"Find" fails, i.e.

If $\forall i = 1, \dots, N$, z_t is α -dept of B_i

$$\Rightarrow \text{for each } i, \|\bar{f} - \underline{f}\|_{B_i} > \alpha \quad (\text{o.w. } \|\bar{f} - \underline{f}\|_{B_i} \leq \alpha \Rightarrow (\bar{f} - \underline{f})^{(z_t)} \leq \alpha)$$

z_t is α -dept of B_i

$$\Rightarrow \|\bar{f} - \underline{f}\|_{S_{t-1}}^2 \geq \sum_{i=1}^N \|\bar{f} - \underline{f}\|_{B_i}^2 > \alpha^2 \cdot N \geq 4\beta$$

Contradiction w/ (*)



Observation 2: Final B_1, \dots, B_N all have size $\leq \text{Edim}(f, \varepsilon)$

Reason: each bin has a sequence of α -indept elts.

Therefore: $\sum_{t=1}^T I(W_{F_t}(z_t) \geq \alpha) \geq \sum_i \text{Final } B_i \text{ size}$

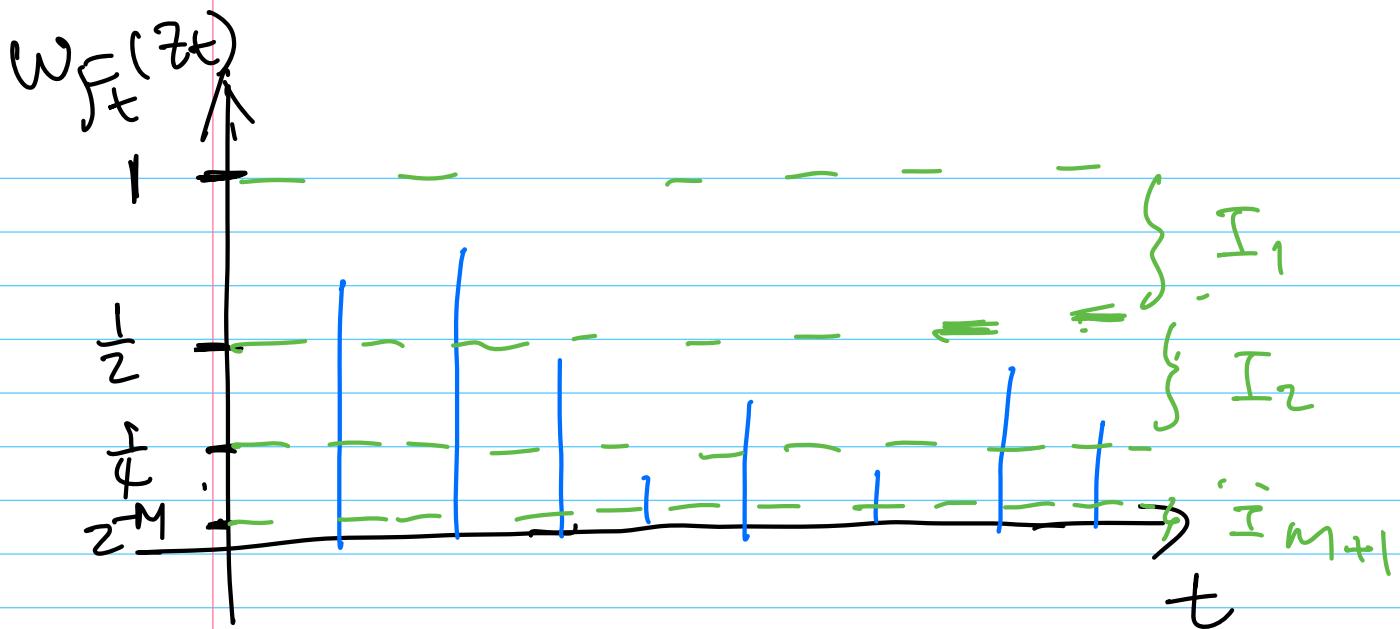
$$\leq \text{Edim}(f, \varepsilon) \lceil \frac{4\beta}{\alpha^2} \rceil$$

Concluding the regret analysis; partition $[T]$ by

$$I_i = \{ t : W_{F_t}(z_t) \in (2^{-i}, 2^{-(i-1)}] \}$$

$$i = 1, 2, \dots, \lfloor \log \frac{1}{\alpha} \rfloor = M \quad (\text{let } \alpha \geq \frac{1}{T} \text{ FBD})$$

$$I_{M+1} = \{ t : W_{F_t}(z_t) \leq 2^{-M} \}$$



"peeling trick".

$$\begin{aligned}
 P\text{Reg}(T) &= \sum_{t \in I_{M+1}} w_{F_t}(z_t) + \sum_{i=1}^M \sum_{t \in I_i} w_{F_t}(z_t) \\
 &\leq |I_{M+1}| \cdot 2^{-M} \quad \underbrace{\sum_{t \in I_i} w_{F_t}(z_t)}_{\leq |I_i| \cdot 2^{-i+1}} \\
 &\leq 2T \alpha . \\
 &\leq \sum_{t=1}^T I(w_{F_t}(z_t) > 2^{-i}) \cdot 2^{-i+1} \\
 &\leq (4\beta z^i + z^{-i}) \text{Edim}(F, \frac{1}{T})
 \end{aligned}$$

$$\leq 2T \alpha + \frac{8\beta \text{Edim}(F, \frac{1}{T})}{\alpha} + \text{Edim}(F, \frac{1}{T}) .$$

$$\text{house } \alpha = \max\left(\frac{1}{T}, 4\sqrt{\frac{\beta \text{Edim}(F, \frac{1}{T})}{T}}\right)$$

$$\leq O\left(\underbrace{\sqrt{\beta E_{dim}(F, \frac{1}{T})} \cdot T}_{\ln|F|} + \underbrace{E_{dim}(F, \frac{1}{T})}_{\text{lower order}}\right)$$

X